# An Open-Source Hierarchical Phrase-Based Machine Translation System

Lane Schwartz
lane@cs.umn.edu

University of Minnesota

Statistical Methods

Hierarchical phrase-based translation



Translation is one of the oldest problems in computer science.

Translation is one of the oldest problems in computer science.

Rule-based translation

Translation is one of the oldest problems in computer science.

Rule-based translation

Requires highly trained bilingual linguists

Translation is one of the oldest problems in computer science.

Rule-based translation

Requires highly trained bilingual linguists

And lots and lots of time for them to develop rules

Translation is one of the oldest problems in computer science.

Rule-based translation

Requires highly trained bilingual linguists

And lots and lots of time for them to develop rules

Hard to extend to new domains

Translation is one of the oldest problems in computer science.

Rule-based translation

Requires highly trained bilingual linguists

And lots and lots of time for them to develop rules

Hard to extend to new domains

Statistical translation

Translation is one of the oldest problems in computer science.

Rule-based translation

Requires highly trained bilingual linguists

And lots and lots of time for them to develop rules

Hard to extend to new domains

► Statistical translation ← We'll use this approach :)

Statistical translation requires large parallel texts.

Corporate documentation

- Corporate documentation
- International news organizations

- Corporate documentation
- International news organizations
- Governments

- Corporate documentation
- International news organizations
- Governments
  - United Nations

- Corporate documentation
- International news organizations
- Governments
  - United Nations
  - Canada

- Corporate documentation
- International news organizations
- Governments
  - United Nations
  - Canada
  - European Union

- Corporate documentation
- International news organizations
- Governments
  - United Nations
  - Canada
  - ► European Union ← We'll use this data :)

Word-by-word translation

- Word-by-word translation
- Word alignments

- Word-by-word translation
- Word alignments
  - Align words by hand

- Word-by-word translation
- Word alignments
  - Align words by hand tedious and time-consuming

- Word-by-word translation
- Word alignments
  - Align words by hand tedious and time-consuming
  - Automatic alignment

- Word-by-word translation
- Word alignments
  - Align words by hand tedious and time-consuming
  - Automatic alignment uses EM on a parallel corpus; see Och & Ney (2000)

- Word-by-word translation
- Word alignments
  - Align words by hand tedious and time-consuming
  - Automatic alignment
    uses EM on a parallel corpus;
    see Och & Ney (2000)
    Many researchers use the freely available GIZA++ tool to
    automatically extract word alignments

- Word-by-word translation
- Word alignments
  - Align words by hand tedious and time-consuming
  - Automatic alignment
    uses EM on a parallel corpus;
    see Och & Ney (2000)
    Many researchers use the freely available GIZA++ tool to
    automatically extract word alignments
- Word alignments can be used in more sophisticated translation models.

Phrases may work better than words

- Phrases may work better than words
- Phrase-based translation

- Phrases may work better than words
- Phrase-based translation
  - ▶ Phrase table

- Phrases may work better than words
- Phrase-based translation
  - Phrase table
  - Reordering model

- Phrases may work better than words
- Phrase-based translation
  - ▶ Phrase table
  - Reordering model
- Phrases can be automatically extracted from word alignments.

► Traditional noisy-channel approach

Traditional noisy-channel approach

$$\arg\max_{e} P(e \mid f) = \arg\max_{e} P(e, f)$$

Traditional noisy-channel approach

$$\arg\max_{e} P(e \mid f) = \arg\max_{e} P(e, f)$$

$$\arg\max_{e} P(e \mid f) = \arg\max_{e} P(e) \times P(f \mid e)$$

Traditional noisy-channel approach

$$\arg\max_{e} \mathsf{P}(e \,|\, f) = \arg\max_{e} \mathsf{P}(e,f)$$
 
$$\arg\max_{e} \mathsf{P}(e \,|\, f) = \arg\max_{e} \mathsf{P}(e) \times \mathsf{P}(f \,|\, e)$$

Log-linear approach

Traditional noisy-channel approach

$$\arg\max_{e} P(e \mid f) = \arg\max_{e} P(e, f)$$

$$\arg\max_{e} P(e \mid f) = \arg\max_{e} P(e) \times P(f \mid e)$$

Log-linear approach

weight 
$$=\prod_i \phi_i^{\lambda_i}$$



# Log-linear features

So what features should our log-linear model use?

### Log-linear features

So what features should our log-linear model use?

- ▶ P<sub>Im</sub>(e)
- ▶ P(f | e)

Features used in noisy channel approach...

### Log-linear features

So what features should our log-linear model use?

- $\triangleright P_{lm}(e)$
- ▶ P(f | e)
- ▶ P(e|f)
- $ightharpoonup P(f_w \mid e_w)$
- $ightharpoonup P(e_w | f_w)$

Features used in noisy channel approach...

...and other features empirically found to be useful!

- Phrase-based translation
  - ▶ Phrase table
  - Reordering model

- Phrase-based translation
  - Phrase table
  - Reordering model
- What if we treat translation as a parsing task?

- Phrase-based translation
  - ▶ Phrase table
  - Reordering model
- What if we treat translation as a parsing task?
  - $\rightarrow$  Phrase table becomes synchronous context free rules

- Phrase-based translation
  - ▶ Phrase table
  - Reordering model
- What if we treat translation as a parsing task?
  - → Phrase table becomes synchronous context free rules
  - → Reordering model becomes implicit in rule applications

Hierarchical phrase-based translation model

Database of synchronous context-free rules

- Database of synchronous context-free rules
- Only two nonterminals!!!

- Database of synchronous context-free rules
- Only two nonterminals!!!
  - X

- Database of synchronous context-free rules
- Only two nonterminals!!!
  - X Used in extracted grammar rules

- Database of synchronous context-free rules
- Only two nonterminals!!!
  - X Used in extracted grammar rules
    - ► S

- Database of synchronous context-free rules
- Only two nonterminals!!!
  - X Used in extracted grammar rules
  - S Allows for serial combination of phrases

▶ Open source implementation of Chiang (2007)

- ▶ Open source implementation of Chiang (2007)
- Implemented in Java

- ▶ Open source implementation of Chiang (2007)
- Implemented in Java
- Designed to be easily extended

- Open source implementation of Chiang (2007)
- Implemented in Java
- Designed to be easily extended
- ▶ Data structures map onto the hypergraph architecture of Huang & Chiang (2005)

- Open source implementation of Chiang (2007)
- Implemented in Java
- Designed to be easily extended
- ▶ Data structures map onto the hypergraph architecture of Huang & Chiang (2005)
  - This allows n-best lists to be easily obtained

- Open source implementation of Chiang (2007)
- Implemented in Java
- Designed to be easily extended
- ▶ Data structures map onto the hypergraph architecture of Huang & Chiang (2005)
  - ▶ This allows n-best lists to be easily obtained
  - N-best lists are needed during parameter tuning

- Open source implementation of Chiang (2007)
- Implemented in Java
- Designed to be easily extended
- Data structures map onto the hypergraph architecture of Huang & Chiang (2005)
  - ▶ This allows n-best lists to be easily obtained
  - N-best lists are needed during parameter tuning
- Uses off-the-shelf minimum error rate trainer for log-linear parameter training

What make this system unique?

► There are two other known implementations of a hierarchical phrase-based system

- ► There are two other known implementations of a hierarchical phrase-based system
- Of the other two...,

- ► There are two other known implementations of a hierarchical phrase-based system
- Of the other two...,
  - ► CMU SAMT Open source, but doesn't implement cube pruning algorithm. C++

- ► There are two other known implementations of a hierarchical phrase-based system
- Of the other two...,
  - ► CMU SAMT Open source, but doesn't implement cube pruning algorithm. C++
  - Chiang's Hiero Closed source, does implement cube pruning algorithm. Python

- ► There are two other known implementations of a hierarchical phrase-based system
- Of the other two...,
  - CMU SAMT Open source, but doesn't implement cube pruning algorithm. C++
  - Chiang's Hiero Closed source, does implement cube pruning algorithm. Python
- Our system...

- ► There are two other known implementations of a hierarchical phrase-based system
- Of the other two...,
  - ► CMU SAMT Open source, but doesn't implement cube pruning algorithm. C++
  - Chiang's Hiero Closed source, does implement cube pruning algorithm. Python
- Our system...
  - Open source

- ► There are two other known implementations of a hierarchical phrase-based system
- Of the other two...,
  - CMU SAMT Open source, but doesn't implement cube pruning algorithm. C++
  - Chiang's Hiero Closed source, does implement cube pruning algorithm. Python
- Our system...
  - Open source
  - ▶ Implements cube pruning algorithm.

- ► There are two other known implementations of a hierarchical phrase-based system
- Of the other two...,
  - CMU SAMT Open source, but doesn't implement cube pruning algorithm. C++
  - Chiang's Hiero Closed source, does implement cube pruning algorithm. Python
- Our system...
  - Open source
  - Implements cube pruning algorithm.
  - Java



Code available through anonymous svn at http://sf.net/projects/nlp-parsers

Questions?

lane@cs.umn.edu